

SkELL: A discovery-based Chinese learning platform

Simon Smith (Coventry University)

Miloš Jakubíček (Lexical Computing)

一

Outline

- Chinese learning
- ICT for Chinese
- What's missing (corpora)
- Planned implementation
- Optimize gramrels (rules) for Chinese
- Implement multi-level segmentation

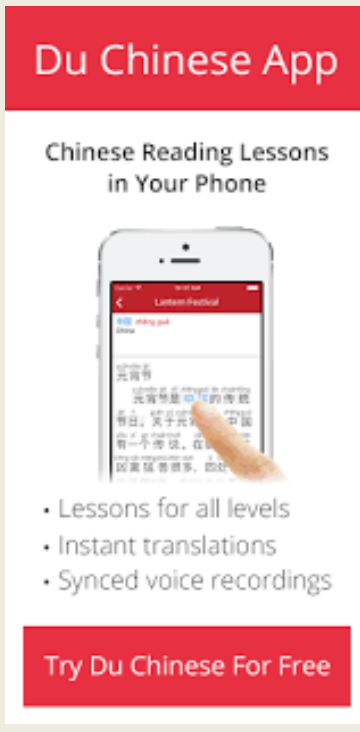
二

Importance of Chinese

- Mandarin Chinese has most native speakers of all languages.
- GCSE entries up 18% in 2015, despite overall decline language exam entries (Guardian 2015).
- Popularity at primary school (Chen 2014)
- Mandarin mastery important for careers (Lo 2016)
- Trade with China
- Confucius Institutes

三

Existing ICT provision



四

Sketch Engine

SkELL

- Acclaimed corpus analysis platform
- English learner dictionaries
- Resources for many languages
- Suitable for expert users
- Can be used by teachers/learners

- Sketch Engine for Language Learning
- Free on web
- English, Czech & Russian currently available
- German & Italian soon
- Suitable for teachers/learners
- Extensible to other languages

References

Chen, T. (2014). Teaching Chinese as a Foreign Language at Primary School in England. *Quarterly Journal of Chinese Studies*, 2(4), 67-83.

Guardian (2015). *GCSE results: fall in numbers taking foreign languages 'a cause for concern'*. Available at <https://www.theguardian.com/education/2015/aug/20/gcse-results-fall-numbers-foreign-languages> [2 January 2017]

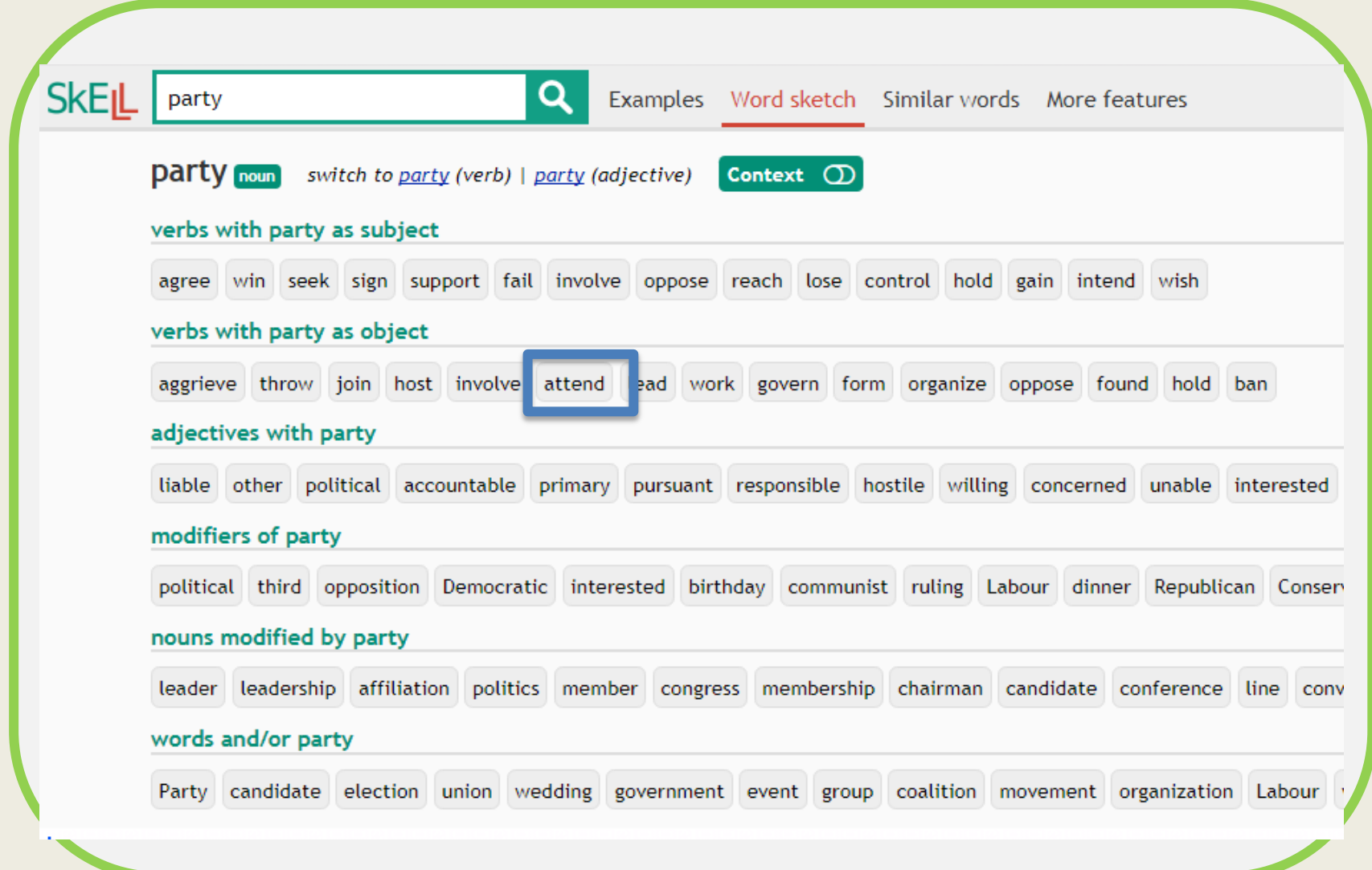
Lo, L. (2016). Challenges faced by Cantonese speakers in a UK university Mandarin course. *Innovative language teaching and learning at university: enhancing participation and collaboration*, 139-145.

Ma, W. Y., & Chen, K. J. (2005). Design of CKIP Chinese word segmentation system. *Chinese and Oriental Languages Information Processing Society*, 14(3), 235-249.

Sun, M., Shen, D., & Tsou, B. K. (1998). Chinese word segmentation without using lexicon and hand-crafted training data. In *Proceedings of the 17th international conference on Computational linguistics-Volume 2* (pp. 1265-1271). Association for Computational Linguistics.

Tseng, H., & Chen, K. J. (2002). Design of Chinese morphological analyzer. In: *Proceedings of the first SIGHAN workshop on Chinese language processing* (pp. 1-7). Association for Computational Linguistics.

First, an English example: studying collocation with SkELL Word sketch



五

Next, a Chinese example: studying collocation of 逮捕 arrest

逮捕 Chinese GigaWord 2 Corpus: Mainland, simplified freq = 8,333 (33.31 per million)

unary rels	Object	Subject	Modifier
Nominalization 58 0.01	4,874 0.58	3,687 0.44	1,643
嫌疑犯 + 110 8.89	警方 + 1,017 9.65	依法 + 551	
犯罪嫌疑人 + 105 8.32	罪名 + 44 8.14	当场 56	
毒贩 45 7.95	警察 + 196 7.47	总共 7	
嫌疑人 43 7.71	罪 29 6.98	予以 29	
肇事者 36 7.62	军警 21 6.77	强行 8	
走私犯 22 7.14	当局 + 143 6.76	予 6	
巴勒斯坦人 89 7.12	犯罪嫌疑人 31 6.73	共 + 113	
头目 32 7.12	巴方 27 6.36	已经 91	
卡拉季奇 22 6.66	预防性 10 6.22	当即 5	
嫌犯 15 6.58	军方 35 6.16	并 81	
凶手 20 6.57	反动派 9 6.12	现已 11	

六

Gramrel examples

Simple VO rule for English

1: "V" "(DET|NUM|ADJ|ADV|N)"* 2:"N"

Object-fronting direct object rule for Chinese

[word="把"|word="将"] NP adv{0,2}
1:"VV" (particle|prep)? NP1 noun

七

Chinese segmentation privileges longer words

八

- Longest match procedure (Tseng & Chen 2002)
- OOV words often long (Ma & Chen 2005)
- Sun et al (1998) use MI segmentation without lexicon
- MI threshold is parametrizable
- Different word lengths could be favoured

[中国人]
[中国][人]
[中][国][人]

Lexicography standard
Learners: word level
Learners: intra-word